

Designing Statistical Estimators That Balance Sample Size, Risk, and Computational Cost

Tongliang Deng

SIST

2018-6-20

Motivation and general method

- ▶ Motivation

When we have a **large amount of data**, we can exploit excess samples to decrease statistical risk, to decrease computational cost, or to trade off between the two.

- ▶ Method

Smooth statistical optimization problems more and more aggressively as the amount of available data increases.

Introduction

- ▶ The Data Model

$$\mathbf{b} = \mathbf{A}\mathbf{x}^\dagger + \mathbf{v}, \mathbf{A} \in \mathbb{R}^{m \times d}$$

Problem: $\hat{\mathbf{x}} := \arg \min_{\mathbf{x}} f(\mathbf{x}), s.t. \|\mathbf{A}\mathbf{x} - \mathbf{b}\| \leq \sqrt{m \cdot R_{\max}} =: \epsilon$

- ▶ Descent cone

The descent cone of a proper convex function

$f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$ at point $\mathbf{x} \in \mathbb{R}^d$ is the convex cone

$$\mathcal{D}(f; \mathbf{x}) := \bigcup_{\tau > 0} \{\mathbf{y} \in \mathbb{R}^d : f(\mathbf{x} + \tau \mathbf{y}) \leq f(\mathbf{x})\}$$

- ▶ Statistical dimension

Let $\mathcal{C} \in \mathbb{R}^d$ be a closed convex cone. Its statistical dimension $\delta(\mathcal{C})$ is

$$\delta(\mathcal{C}) := \mathbb{E}_{\mathbf{g}}[\|\Pi_{\mathcal{C}}(\mathbf{g})\|^2]$$

Introduction

- Phase Transition

Whenever $m < \delta$,

$$\lim_{\sigma \rightarrow 0} \frac{\mathbb{E}_v[R(x^*)|\mathbf{A}]}{\sigma^2} = 0$$

with probability p_1 .

Whenever $m > \delta$,

$$\left| \lim_{\sigma \rightarrow 0} \frac{\mathbb{E}_v[R(x^*)|\mathbf{A}]}{\sigma^2} - \left(1 - \frac{\delta}{m}\right) \right| \leq tm^{-1}\sqrt{d}$$

with probability p_2 .

Geometric opportunity

Enlarging convex constraint sets can make corresponding statistical optimization problems easier to solve. These geometric deformations, however, create a loss of statistical accuracy.

- By enlarging the sublevel sets of the regularizer f , we increase the statistical dimension of the descent cone of f at \mathbf{x}^h .

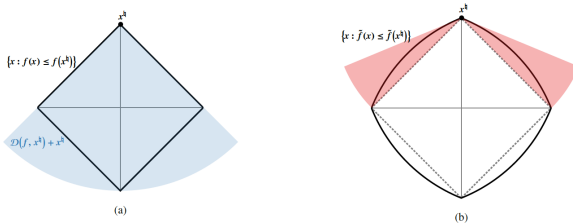


Fig. 1. **A geometric opportunity.** Panel (a) illustrates the sublevel set and descent cone of a regularizer f at the point \mathbf{x}^h . Panel (b) shows a relaxed regularizer \tilde{f} with larger sublevel sets. The shaded area indicates the difference between the descent cones of \tilde{f} and f at \mathbf{x}^h . Fact III.3 shows how this difference in the size of the descent cones translates into a difference in statistical accuracy. We may compensate for this loss of statistical accuracy by choosing a relaxation \tilde{f} that allows us to solve the optimization problem faster.

Computational opportunity

- Primal problem

$$\hat{x}_\mu := \arg \min_x f_\mu(x), \text{ subject to } \|\mathbf{Ax} - \mathbf{b}\| \leq \epsilon$$

- Dual problem

$$\text{maximize } g_\mu(z) := \inf_x \{f_\mu(\mathbf{x}) - \langle z, \mathbf{Ax} - \mathbf{b} \rangle - \epsilon \|z\|\},$$

where z is the dual variable.

- Rewrite the dual problem

$$\begin{aligned} g_\mu(z) &= \inf_x \{f_\mu(\mathbf{x}) - \langle \mathbf{A}^T z, \mathbf{x} \rangle\} + \langle z, \mathbf{b} \rangle - \epsilon \|z\| \\ &= \underbrace{-f_\mu^*(\mathbf{A}^T z)}_{\tilde{g}_\mu(z)} + \underbrace{\langle z, \mathbf{b} \rangle - \epsilon \|z\|}_{h(z)}, \end{aligned}$$

where f_μ^* is the convex conjugate of f_μ . \tilde{g}_μ has a Lipschitz gradient $L_\mu = \mu^{-1} \|\mathbf{A}\|^2$.

Auslender-TeBoulle Algorithm

Algorithm 1. Auslender–Teboulle

Input: measurement matrix \mathbf{A} , observed vector \mathbf{b} , parameter ϵ

```
1:  $\mathbf{z}_0 \leftarrow \mathbf{0}, \bar{\mathbf{z}}_0 \leftarrow \mathbf{z}_0, \theta_0 \leftarrow 1$   
2: for  $k = 0, 1, 2, \dots$  do  
3:    $\mathbf{y}_k \leftarrow (1 - \theta_k)\mathbf{z}_k + \theta_k \bar{\mathbf{z}}_k$   
4:    $\mathbf{x}_k \leftarrow \arg \min_{\mathbf{x}} \{f(\mathbf{x}) + \langle \mathbf{y}_k, \mathbf{b} - \mathbf{A}\mathbf{x} \rangle\}$   
5:    $\bar{\mathbf{z}}_{k+1} \leftarrow \text{Shrink}(\bar{\mathbf{z}}_k - (\mathbf{b} - \mathbf{A}\mathbf{x}_k)/(L_\mu \cdot \theta), \epsilon/(L_\mu \cdot \theta))$   
6:    $\mathbf{z}_{k+1} \leftarrow (1 - \theta_k)\mathbf{z}_k + \theta_k \bar{\mathbf{z}}_{k+1}$   
7:    $\theta_{k+1} \leftarrow 2/(1 + (1 + 4/\theta_k^2)^{1/2})$   
8: end for
```

There is an bug at line 5, it should be

$$\bar{\mathbf{z}}_{k+1} \leftarrow \text{Shrink}(\bar{\mathbf{z}}_k + (\mathbf{b} - \mathbf{A}\mathbf{x}_k)/(L_\mu \cdot \theta), \epsilon/(L_\mu \cdot \theta))$$

Time-date tradeoff

- ▶ A dual-smoothing method

Given a regularizer f in the problem, introduce a family $\{f_\mu : \mu > 0\}$ of strong convex majorants:

$$f_\mu(\mathbf{x}) := f(\mathbf{x}) + \frac{\mu}{2} \|\mathbf{x}\|^2$$

These majorants have larger sublevel sets, and their descent cones have larger statistical dimension.

- ▶ Choosing a smoothing parameter

1. Constant smoothing, choose a constant value of μ .

2. Constant risk, $\frac{\delta(\mathcal{D}(f_\mu; \mathbf{x}^\natural))}{m} = \frac{\bar{\delta}}{\bar{m}}$.

3. A tunable balance, $\frac{\delta(\mathcal{D}(f_\mu; \mathbf{x}^\natural))}{m} = \frac{\bar{\delta}}{\bar{m} + (m - \bar{m})^\alpha}$.

Experimental setup

- ▶ Intel Core 5, memory 8GB, matlab2016Rb, Linux.
- ▶ Sparse vector regression.

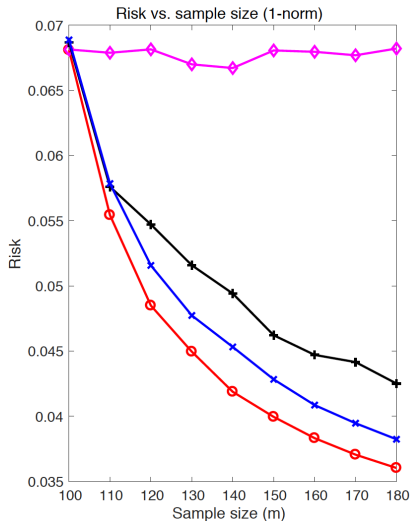
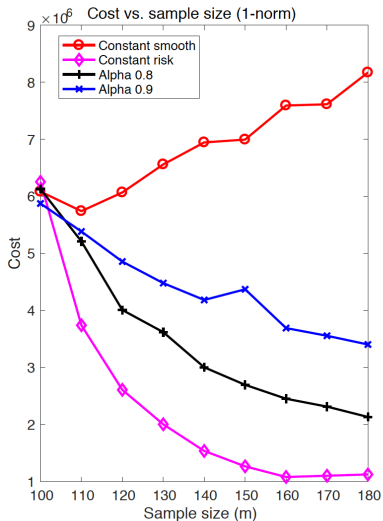
$$\delta(\mathcal{D}(f_\mu; \mathbf{x})) = \mathbf{d} \cdot \psi(\rho)$$

where $\psi: [0, 1] \rightarrow \mathbb{R}$ is the function given by

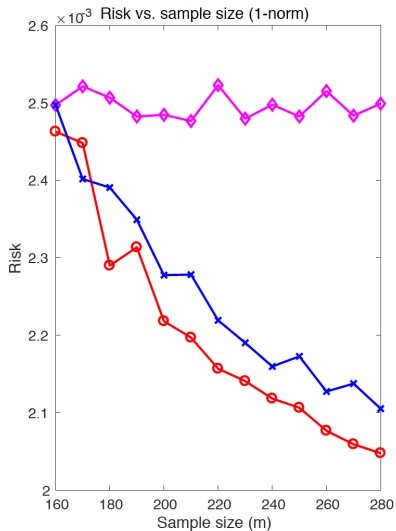
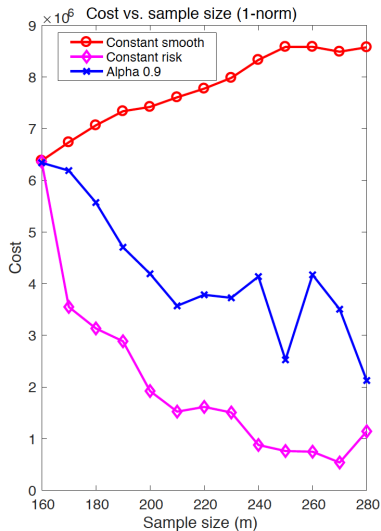
$$\psi(\rho) = \inf_{\tau \geq 0} \left\{ \rho \left[1 + \tau^2 (1 + \mu \|\mathbf{x}\|_{\ell_\infty})^2 \right] \right. \\ \left. + (1 - \rho) \sqrt{\frac{2}{\pi}} \int_{\tau}^{\infty} (u - \tau)^2 e^{-u^2/2} du \right\}.$$

- ▶ Low-rank matrix regression.
Refer to the paper

Result of sparse vector regression



Result of Low-rank matrix regression



Thank you!